

Translational Medicine Research

Series Editors: Zhu Chen · Xiaoming Shen  
Saijuan Chen · Kerong Dai



# Advance in Structural Bioinformatics

结构生物信息学 (英文版)

Dongqing Wei · Qin Xu · Tangzhen Zhao · Hao Dai *Editors*

魏冬青 徐 沁 赵唐祯 戴 昊 编著



上海交通大学出版社  
SHANGHAI JIAO TONG UNIVERSITY PRESS



Springer

# Advance in Structural Bioinformatics

结构生物信息学  
(英文版)

Dongqing Wei · Qin Xu · Tangzhen Zhao · Hao Dai *Editors*

魏冬青 徐沁 赵唐祯 戴昊 编著



上海交通大学出版社  
SHANGHAI JIAO TONG UNIVERSITY PRESS



Springer

图书在版编目 (C I P) 数据

结构生物信息学 : 英文 / 魏冬青等编著. -- 上海 : 上海交通大学出版社, 2015

(转化医学出版工程 / 陈竺, 沈晓明, 陈赛娟, 戴尅戎主编)

ISBN 978-7-313-13982-5

I. ①结… II. ①魏… III. ①生物结构—生物信息论—研究—英文 IV. ①Q617②Q811.4

中国版本图书馆CIP数据核字 (2015) 第253419号

Not for sale outside the Mainland of China

(Not for sale in Hong Kong SAR, Macau SAR, and Taiwan,  
and all countries except the Mainland of China)

结构生物信息学 (英文版)

编 著: 魏冬青 徐沁 赵唐祯 戴昊

出版发行: 上海交通大学出版社

邮政编码: 200030

出 版 人: 韩建民

印 制: 上海锦佳印刷有限公司

开 本: 787mm × 1092mm 1/16

字 数: 762千字

版 次: 2015年12月第1版

书 号: ISBN 978-7-313-13982-5/Q

定 价: 248.00元

地 址: 上海市番禺路951号

电 话: 021-64071208

经 销: 全国新华书店

印 张: 24.75

印 次: 2015年12月第1次印刷

版权所有 侵权必究

告读者: 如发现本书有印装质量问题请与印刷厂质量科联系

联系电话: 021-56401314

# **Advances in Experimental Medicine and Biology**

**Volume 827**

## **Series editors**

**Irwin R. Cohen, Rehovot, Israel**

**N.S. Abel Lajtha, Orangeburg, USA**

**Rodolfo Paoletti, Milan, Italy**

**John D. Lambris, Philadelphia, USA**

More information about this series at <http://www.springer.com/series/5584>

## Foreword

Structural bioinformatics, one of the hot spots of bioinformatics, is experiencing a rapid development in recent years. In the genome era, proteomics, genomics, and other data increase dramatically, providing a basis to clarify the problem of essential physiological functions of nucleic acids, proteins, and other biological macromolecules. Relative to the traditional sequence-based bioinformatics, structural bioinformatics focuses mainly on the exploration of the structure and function of biological macromolecules and their dynamic properties. Many human serious diseases are generally associated with some of the key enzymes, ion channels, or associated regulatory proteins. So, most of the new drug research is designed targeting on these proteins. Compared to the previous experimental approaches and sequence analysis, a more comprehensive knowledge of the physiological and pathological mechanism of the drug and the target protein could be obtained from the view of the spatial three-dimensional structure of these molecules and their dynamic structural changes.

The primary problem structural bioinformatics has been trying to solve is that we can build a protein model to fully reveal the nature of its structure and function through the extraction and analysis of the current high-throughput data of biological macromolecules, combining with structural biology knowledge and bioinformatics methods. Besides, to deduce and predict the unknown molecular structure and function based on the known one, and further to realize computer-aided the design and customization of the structure of protein complexes is a long-term goal.

This book represents comprehensive introduction and latest progresses in various aspects of structural bioinformatics. It covers not only the knowledge of mathematical and physical modeling theory, but also the computational methods and its applications in structural bioinformatics. More important, it takes the latest research achievements from the leading groups in this field as examples to illustrate the basic molecular dynamic theory. The content of this book mainly includes the basic knowledge of structural bioinformatics, genomics and proteomics sequence acquiring and analysis, structures of protein, DNA and RNA, basic methods of

molecular dynamic simulations and conformation search, the application examples of computing simulation methods and the structure-based drug design, recent research progress, and future prospects.

We are most grateful to professors and students in the class of “Structural Bioinformatics” at Shanghai Jiao Tong University, where the main contents of this book are accumulated.

Minhang, Shanghai, January 2014

Dongqing Wei

# Acknowledgments

In the process of putting this book together, we are much indebted to many people who gave generous support. I would like to express deepest gratitude to the many friends who saw me through this book; to all those who provided support, talked things over, read, wrote, offered comments, allowed me to quote their remarks and assisted in the editing, proofreading and design.

I would like to especially thank the following authors, who were invited to contribute some chapters to this book. They are all from leading research groups in the field of structural bioinformatics in the world, with some of whom I have had the honor to collaborate, i.e., authors of the following chapters

Chapter 2 JVM: Java Visual Mapping Tool for Next Generation Sequencing Read. Ye Yang, Juan Liu\*.

Chapter 3 Advancement of Polarizable Force Field and Its Use for Molecular Modeling and Design. Peijun Xu, Huiying Chu, Beibei Li, Yingchen Mao, Yang Ding, Guohui Li\*.

Chapter 4 Systematic Methods for Defining Coarse-Grained Maps in Large Biomolecules. Zhiyong Zhang\*.

Chapter 5 Quantum Calculation of Protein NMR Chemical Shifts Based on the Automated Fragmentation Method. Tong Zhu, John Z.H. Zhang, Xiao He\*.

Chapter 7 Extended Structure of Rat Islet Amyloid Polypeptide in Solution. Lei Wei, Ping Jiang, Malathy Sony Subramanian Manimekalai, Cornelia Hunke, Gerhard Grüber, Konstantin Pervushin, Yuguang Mu\*.

Chapter 8 Folding Mechanisms of Trefoil Knot Proteins Studied by Molecular Dynamics Simulations and Go-model. Xue Wu, Ting Fu, Zhilong Xiu, Guohui Li\*.

Chapter 9 Binding Induced Intrinsically Disordered Protein Folding with Molecular Dynamics Simulation. Haifeng Chen\*.

Chapter 10 Theoretical Studies on the Folding Mechanisms for Different DNA G-quadruplexes. Xue Wu, Ting Fu, Hujun Shen, Zhilong Xiu, Guohui Li\*.

Chapter 11 RNA Folding: Structure Prediction, Folding Kinetics and Ion Electrostatics. Zhijie Tan\*, Wenbing Zhang\*, Yazhou Shi, Fenghua Wang.



Chapter 12 Binding Modes and Interaction Mechanism Between Different Base Pairs and Methylene Blue Trihydrate: A Quantum Mechanics Study. Huiying Chu, Jinguang Wang, Yong Xu, Hujun Shen, Guohui Li\*.

Chapter 15 Evolutionary Optimization of Transcription Factor Binding Motif Detection. Zhao Zhang, Ze Wang, Guoqin Mai, Youxi Luo, Miaomiao Zhao, Fengfeng Zhou\*.

Chapter 16 Prediction of Serine/Threonine Phosphorylation Sites in Bacteria Proteins. Zhengpeng Li, Ping Wu, Yuanyuan Zhao, Zexian Liu\*, Wei Zhao\*.

Chapter 22 Bayesian Analysis of Complex Interacting Mutations in HIV Drug Resistance and Cross-Resistance. Ivan Kozyryev, Jing Zhang\*.

I would like to thank Springer and Shanghai Jiao Tong University Press for providing me with the opportunity to edit this book.

# Contents

<b>1</b>	<b>Introduction to Structural Bioinformatics</b> . . . . .	<b>1</b>
	Qin Xu, Hao Dai, Tangzhen Zhao and Dongqing Wei	

## **Part I Advances in Methods for Structural Bioinformatics**

<b>2</b>	<b>JVM: Java Visual Mapping Tool for Next Generation Sequencing Read</b> . . . . .	<b>11</b>
	Ye Yang and Juan Liu	
<b>3</b>	<b>Advancement of Polarizable Force Field and Its Use for Molecular Modeling and Design</b> . . . . .	<b>19</b>
	Peijun Xu, Jinguang Wang, Yong Xu, Huiying Chu, Jiahui Liu, Meixia Zhao, Depeng Zhang, Yingchen Mao, Beibei Li, Yang Ding and Guohui Li	
<b>4</b>	<b>Systematic Methods for Defining Coarse-Grained Maps in Large Biomolecules</b> . . . . .	<b>33</b>
	Zhiyong Zhang	
<b>5</b>	<b>Quantum Calculation of Protein NMR Chemical Shifts Based on the Automated Fragmentation Method</b> . . . . .	<b>49</b>
	Tong Zhu, John Z.H. Zhang and Xiao He	
<b>6</b>	<b>Applications of Rare Event Dynamics on the Free Energy Calculations for Membrane Protein Systems</b> . . . . .	<b>71</b>
	Yukun Wang, Ruoxu Gu, Huaimeng Fan, Jakob Ulmschneider and Dongqing Wei	

## Part II 3D-Structure Prediction and Folding Mechanism of Biological Macromolecules

- 7 **Extended Structure of Rat Islet Amyloid Polypeptide in Solution** ..... 85  
Lei Wei, Ping Jiang, Malathy Sony Subramanian Manimekalai, Cornelia Hunke, Gerhard Grüber, Konstantin Pervushin and Yuguang Mu
- 8 **Folding Mechanisms of Trefoil Knot Proteins Studied by Molecular Dynamics Simulations and Go-model** ..... 93  
Xue Wu, Peijun Xu, Jinguang Wang, Yong Xu, Ting Fu, Depeng Zhang, Meixia Zhao, Jiahui Liu, Hujun Shen, Zhilong Xiu and Guohui Li
- 9 **Binding Induced Intrinsically Disordered Protein Folding with Molecular Dynamics Simulation** ..... 111  
Haifeng Chen
- 10 **Theoretical Studies on the Folding Mechanisms for Different DNA G-quadruplexes** ..... 123  
Xue Wu, Peijun Xu, Jinguang Wang, Yong Xu, Ting Fu, Meixia Zhao, Depeng Zhang, Jiahui Liu, Hujun Shen, Zhilong Xiu and Guohui Li
- 11 **RNA Folding: Structure Prediction, Folding Kinetics and Ion Electrostatics** ..... 143  
Zhijie Tan, Wenbing Zhang, Yazhou Shi and Fenghua Wang

## Part III The Interactions Between Biological Macromolecules and Ligands

- 12 **Binding Modes and Interaction Mechanism Between Different Base Pairs and Methylene Blue Trihydrate: A Quantum Mechanics Study** ..... 187  
Peijun Xu, Jinguang Wang, Yong Xu, Huiying Chu, Hujun Shen, Depeng Zhang, Meixia Zhao, Jiahui Liu and Guohui Li
- 13 **Drug Inhibition and Proton Conduction Mechanisms of the Influenza A M2 Proton Channel** ..... 205  
Ruoxu Gu, Limin Angela Liu and Dongqing Wei

- 14 Exploring the Ligand-Protein Networks in Traditional Chinese Medicine: Current Databases, Methods and Applications** 227  
Mingzhu Zhao and Dongqing Wei

#### **Part IV Functional Analysis of Biological Macromolecules**

- 15 Evolutionary Optimization of Transcription Factor Binding Motif Detection** . . . . . 261  
Zhao Zhang, Ze Wang, Guoqin Mai, Youxi Luo,  
Miaomiao Zhao and Fengfeng Zhou
- 16 Prediction of Serine/Threonine Phosphorylation Sites in Bacteria Proteins** . . . . . 275  
Zhengpeng Li, Ping Wu, Yuanyuan Zhao,  
Zexian Liu and Wei Zhao
- 17 Bioinformatics Tools for Discovery and Functional Analysis of Single Nucleotide Polymorphisms** . . . . . 287  
Li Li and Dongqing Wei
- 18 An Application of QM/MM Simulation: The Second Protonation of Cytochrome P450** . . . . . 311  
Peng Lian and Dongqing Wei

#### **Part V Application of Structural Bioinformatics in Drug Design**

- 19 Recent Progress on Structural Bioinformatics Research of Cytochrome P450 and Its Impact on Drug Discovery** . . . . . 327  
Tao Zhang and Dongqing Wei
- 20 Human Cytochrome P450 and Personalized Medicine** . . . . . 341  
Qi Chen and Dongqing Wei
- 21 The  $\alpha 7$  nAChR Selective Agonists as Drug Candidates for Alzheimer's Disease** . . . . . 353  
Huaimeng Fan, Ruoxu Gu and Dongqing Wei
- 22 Bayesian Analysis of Complex Interacting Mutations in HIV Drug Resistance and Cross-Resistance** . . . . . 367  
Ivan Kozyryev and Jing Zhang

# Chapter 1

## Introduction to Structural Bioinformatics

Qin Xu, Hao Dai, Tangzhen Zhao and Dongqing Wei

**Abstract** Structural Bioinformatics is one of the hot spots of interdisciplinary sciences and obtained amazing advances in recent years. The first chapter overviews the concept of structural bioinformatics, and briefly describe the contents of this book. The interdisciplinary corporations make it difficult to further divide structural bioinformatics, so the chapters in this book are roughly separated according to the different fields of their applications. That is, fundamental developments in methods of structural bioinformatics, tertiary structure prediction and folding mechanism analysis, the binding mechanism and the interactions between biological macromolecules and ligands, structure-based functional analysis of biological macromolecules, as well as the applications in drug design.

**Keywords** Structural bioinformatics • Structure of macromolecules • Structure-based drug design

### 1.1 What Is Structural Bioinformatics

Structural Bioinformatics is generally looked as a branch of bioinformatics mainly about problems of structural biology, which the word “structural” is referred to here. In the early days, it was also named as “computational structural biology”, using the distinctive techniques of computational molecular simulations. And the

---

Q. Xu · H. Dai · T. Zhao · D. Wei (✉)

State Key Laboratory of Microbial Metabolism, College of Life Sciences  
and Biotechnology, Shanghai Jiao Tong University, Shanghai, China  
e-mail: dqwei@sjtu.edu.cn

research interests were mainly focused in analysis and prediction of the three-dimensional structures and related functions of biological macromolecules such as proteins, RNA, and DNA.

However, the fast developments in technologies and combinations with other fields make structural bioinformatics more and more diverse and interdisciplinary. Mathematics, statistics, informational sciences, bioinformatics, biophysics, computational chemistry, structural biology, enzymology, medical engineering, pharmaceutical sciences, and much more other disciplines are making contributions to structural bioinformatics. In the meanwhile, its applications are expanding into much more fields, like comparisons of overall folds and local motifs of both primary, secondary and tertiary structures, structural and functional predictions, molecular mechanism of folding/unfolding of macromolecules, evolution and bioengineering, binding interactions in the macromolecules complexes like drug-target complex, molecular mechanism of enzymatic catalysis, as well as other structure-function relationships. In addition to its wide application in the researches of biological sciences, it is showing more power in the industries of bioengineering and drug developments.

The award of 2013 Nobel Prize in Chemistry to Martin Karplus, Michael Levitt, and Arieh Warshel “for the development of multiscale models for complex chemical systems” is, in a way, recognition of the importance of computational techniques in chemistry and biology. However, the computational methods are not opposite to the experimental ones, but complimentary and embedded into them, boosting the developments of more new and advanced techniques and methods to be used. Finally, these new methods might result into a new field of technologies or sciences. Here, structural bioinformatics is a successful example: the advances in this interdisciplinary science have gradually made it an unignorable discipline.

## 1.2 What Is in This Book

The fast developments in structural bioinformatics attracted more research interests, brought more collaboration from different scientific scopes, and resulted into more advances, both in methodology and in applications. In this book, some of these new advances in structural bioinformatics are introduced, so that the researcher interested in this new field could get some new idea in the scientific developments or interdisciplinary collaborations from these successful examples.

The diverse interdisciplinary combinations make it difficult to trace the development of structural bioinformatics in a single line or divide it into sub-disciplines. But the emergence of structural bioinformatics could be somewhat simply explained as the application of new bioinformatic technologies into the research of structural biology. Therefore, in this book the chapters are organized roughly according to the different applications of the new techniques, additional to those advances with more emphasis on methodology, which are described briefly in the sections below.

### 1.2.1 Part I: Advances in Methods for Structural Bioinformatics

In Part I, we first introduced several new advances to improve the methodology of structural bioinformatics in different fields, like sequencing, molecular simulation and *in silico* computational chemistry.

Chapter 2 is about program JVM, a powerful tool for mapping next generation sequencing read to reference sequence. It can deal with millions of short read generated by sequence alignment using the Illumina sequencing technology, employing seed index strategy and octal encoding operations for sequence alignments. It is implemented in Java and designed as a desktop application, which supports reads capacity from 1MB to 10 GB. JVM is useful for DNA-Seq, RNA-Seq when dealing with single-end resequencing.

Molecular simulation is always one of the major methods of structural bioinformatics. The contribution of molecular simulation to the developments of chemistry was recently recognized by the 2013 Nobel Prize in Chemistry. The various methods of simulations have covered a diversity of biological scales now. The most popular method, the classical molecular dynamics is fully depended on the force field used. One of the current hot spots of force fields is how to deal with the influence of the electrostatic polarization. In Chap. 3, we review the history of the classical force fields and polarizable force fields, together with its application on small molecules and biological macromolecules simulation, as well as molecular design. In the meantime, various coarse-grained (CG) approaches have also attracted rapidly growing interest in this field of research, because they enable simulations of large biomolecules over longer effective timescales than all-atom molecular dynamics (MD) simulations. Chapter 4 reviews the recent development of a novel and systematic method for constructing CG representations of arbitrary biomolecules, which preserves large-scale and functionally relevant essential dynamics (ED) at the CG level. This method may serve as a very useful tool for the identification of functional dynamics of large biomolecules at the CG level. In Chap. 5, techniques of rare event dynamics are reviewed, followed by further discussion on the intrinsic difficulties to calculate free energy of rare events and the introduction of several well-developed free energy calculation methods. Then several examples of free energy calculations are illustrated, like the calculations on the drug binding in the M2 proton channel, as well as the insertion and association of membrane proteins and membrane active peptides.

In Chap. 6, the automatic fragmentation quantum mechanics/molecular mechanics (AF-QM/MM) is introduced to calculate the *ab initio* NMR chemical shifts so as to improve protein structure determination and refinement. Using the Poisson-Boltzmann (PB) model and first solvation water molecules, the influence of solvent effect is also discussed. Benefit from the fragmentation algorithm, the AF-QM/MM approach is computationally efficient, linear-scaling with a low pre-factor, and massively parallel.

### ***1.2.2 Part II: 3D-Structure Prediction and Folding Mechanism of Biological Macromolecules***

Part II focuses on one of the main applications of structural bioinformatics since its early days, that is, the structural prediction and analysis on the mechanisms of folding/unfolding of biological macromolecules. Without good understanding of the structure of the research objects, any in-depth study is questionable.

The case in Chap. 7 is about the research of the extend structure of human islet amyloid polypeptide (hIAPP). The human IAPP aggregates easily, so it is difficult to characterize its structural features by standard biophysical tools. The problem was solved by using rat version of IAPP (rIAPP) as substitute which differs from human IAPP by six amino acids and is not prone to aggregation and does not form amyloid fibrils and similar to human IAPP, it demonstrates random-coiled nature. However, the overall shape of it in solution still remains elusive. Using small angle X-ray scattering (SAXS) measurements combined with nuclear magnetic resonance (NMR) and molecular dynamics simulations (MD) the solution structure of rIAPP was studied and an overall random-coiled feature with residual helical propensity in the N-terminus was confirmed eventually.

The application of structural bioinformatics on the analysis of protein folding mechanisms is illustrated by two examples in Chaps. 8 and 9. In Chap. 8, the folding mechanism of two trefoil knot proteins was simulated under high temperature using all-atom Gō-model. Similar results of the folding process were obtained for the two proteins. That is, the contacts in  $\beta$ -sheet are important to the formation of knot protein. Without these contacts, the knot protein would be easy to untie. In Chap. 9, the folding mechanism of intrinsically disordered proteins upon partner binding was simulated under room temperature as well as high-temperature. The former suggests both nonspecific and specific interactions between the intrinsically disordered proteins and the partner, while the latter shows the kinetics of a two-state process for both the unfolding of apo-states and the unbinding of the bound states. Based on the results of the unfolding processes, the folding pathway of bound intrinsically disordered protein was proposed as: unfolded state, secondary structure folding, tertiary folding, partner binding, and finally to the folded state. In addition, induced-fit mechanism was suggested for the specific recognition between intrinsically disordered protein and its partner using Kolmogorov-Smirnov (KS)  $P$  test analysis.

In the rest part of Part II, we presented applications of structural bioinformatics in the studies of DNA and RNA folding. Chapter 10 discusses the folding mechanisms of different DNA G-quadruplexes, which could be a promising anticancer target. In this study, the folding of the thrombin aptamer, Form1 and Form3 G-quadruplexes were simulated with all-atom Gō-model and analyzed by the energy landscape theory, and all were suggested to be a two-state mechanism: the compact structures are formed in the initial stage of the folding process, then they are folded to the native states through the formation of G-triplex structures. The free energy barrier to fold Form 3 G-quadruplex is higher than those to fold thrombin aptamer and Form1, suggesting higher stability of Form 3 G-quadruplex



than those of the other two G-quadruplexes. In Chap. 11, we review the recent experimental and theoretical progress, especially the theoretical modeling of the three major problems in RNA folding: structure prediction, folding kinetics and influence of ion electrostatics.

### ***1.2.3 Part III: The Interactions Between Biological Macromolecules and Ligands***

Part III emphasizes on the interactions between macromolecules like protein or DNA/RNA and small ligand molecules, especially possible drug like compounds.

Chapter 12 studies the interactions between DNA base pairs and methylene blue trihydrate, a dye and therapeutic agent possibly to be inserted into two adjacent DNA base pairs. Thus it is called a DNA intercalator. Its binding mode with different base pairs was evaluated and compared using a series of quantum mechanical methods, including various semi-empirical methods, DFT methods and *ab initio* methods. The results showed that the DFT method WB97XD with 6-311+G\* basis set best reproduced the result of the expensive *ab initio* method MP2 and determined that the best binding mode was into the AA-TT base pair according to the binding energies and charge density analyses.

Chapter 13 is about the influenza A virus matrix protein 2 (M2 protein), a pH-regulated proton channel crucial to the viral infection and replication. In this chapter, the experimental and computational studies of the two possible drug binding sites on the M2 protein were reviewed to explain the mechanisms for inhibitors to prevent proton conduction, the recent molecular dynamics simulations of the interactions between amantadine and drug-resistant mutant channels were summarized to propose mechanisms for drug resistance, and two proton conduction mechanisms in debate were discussed to further illustrate the applications of structural bioinformatics to understand the structure and functions of this interesting membrane protein.

In Chap. 14, the studies of protein-ligand interactions are in a totally different way, in which massive information about ligand bioactivity and the target protein structures were summarized into the ligand-protein networks so as to elucidate possible “multi-component—multi-target” mechanism of the traditional Chinese medicine (TCM) from its complex composition and unclear pharmacology.

### ***1.2.4 Part IV: Functional Analysis of Biological Macromolecules***

It is generally believed that the functions of biological macromolecules are in some ways determined by their structures, including primary structures, secondary structures and tertiary structures. Therefore, one of the major applications of structural bioinformatics is to analyze or predict the functions, activities,